Formal Methods for XML: Algorithms & Complexity

S. Margherita di Pula August 2004

Thomas Schwentick

XML

Example Document

```
(Composer)
```

```
(Name) Claude Debussy (/Name)
(Vita)
   (Born) (When) August 22, 1862 (/When)(Where) Paris (/Where)(/Born)
   (Married)(When) October 1899 (/When)(Whom) Rosalie(/Whom)(/Married)
   (Married)(When) January 1908 (/When)(Whom) Emma (/Whom)(/Married)
   (Died)(When) March 25, 1918 (/When)(Where) Paris (/Where) (/Died)
(/Vita)
(Piece)
   (PTitle) La Mer (/PTitle)
   (PYear) 1905 (/PYear)
   (Instruments) Large orchestra (/Instruments)
   (Movements) 3 (/Movements)
```

```
...
⟨/Piece⟩
```

```
(/Composer)
```

. . .

XML





Four important kinds of XML processing Validation

Check whether an XML document is of a given type

Navigation

Select a set of positions in an XML document

Querying

Extract information from an XML document

Transformation

Construct a new XML document from a given one

XML Processing _____

Four important kinds of XML processing	and	their	languages
Validation	DTD,	XML	Schema
Check whether an XML document is of a given	type		
Navigation			XPath
Select a set of positions in an XML document			
Querying			XQuery
Extract information from an XML document			
Transformation			XSLT
Construct a new XML document from a given of	one		

Example document

```
(Composer)
   (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     (Movements) 3 (/Movements)
     . . .
  \langle / Piece \rangle
(/Composer)
```

Schwentick

Validation: DTD

DTD

DTDs describe types of

Example document

XML documents

```
(Composer)
  (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     (Movements) 3 (/Movements)
     . . .
  \langle / Piece \rangle
(/Composer)
```

Validation: DTD

DTD

DTDs describe types of

XML documents

(Composer)

(Piece)

Example document

Example

```
<
```

```
<!DOCTYPE Composers [

<!ELEMENT Composers (Composer*)>

<!ELEMENT Composer (Name, Vita, Piece*)>

<!ELEMENT Vita (Born, Married*, Died?)>

<!ELEMENT Born (When, Where)>

<!ELEMENT Born (When, Where)>

<!ELEMENT Died (When, Where)>

<!ELEMENT Piece (PTitle, PYear,

Instruments, Movements)>

]>
```

Navigation: XPath

Example document

```
(Composer)
  (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     (Movements) 3 (/Movements)
     . . .
  \langle / Piece \rangle
(/Composer)
```

Schwentick

Navigation: XPath _____

	XPath		
	XPath expressions select sets of	nodes of	
Example doc	XML documents by specifying na	avigational	
Composer> (Name) Claude Debussy (/Name)	patterns	0	
(Vita)			
(Born) (When) August 22, 1862 (/Wh	en〉〈Where〉 <mark>Paris</mark> 〈/Where〉〈/Born〉		
(Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)			
(Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)			
(Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)			
$\langle /Vita \rangle$			
<pre>〈Piece〉</pre>			
<pre>〈PTitle〉La Mer </pre>			
(PYear) 1905 (/PYear)			
(Instruments) Large orchestra (/Instruments)			
<pre>〈Movements〉 3 〈/Movements〉</pre>			
/Composer>			

Navigation: XPath ____

	XPath	
	XPath expressions select sets of	of nodes of
Example doc	XML documents by specifying r	navigational
Composer> 〈Name〉 Claude Debussy 〈/Name〉	patterns	-
<pre>{Vita}</pre>	ien〉〈Where〉 Paris 〈/Where〉〈/Born〉 n〉〈Whom〉 Rosalie〈/Whom〉〈/Married〉 n〉〈Whom〉 Emma 〈/Whom〉〈/Married〉 〉〈Where〉 Paris 〈/Where〉 〈/Died〉	>
<pre>{Instruments> Large orchestra {/Instrum {Movements> 3 {/Movements> {/Piece></pre>	ients>	
/Composer〉		
•	Exampl	e query
	//Vita	/Died/*

Schwentick

Navigation: XPath



Example document

```
(Composer)
   (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     \langle Movements \rangle 3 \langle Movements \rangle
  \langle / Piece \rangle
(/Composer)
```

Schwentick

XQuery

XQuery is a full-fledged XML query language

```
\langle Composer \rangle
```

```
(Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     \langle Movements \rangle 3 \langle Movements \rangle
  \langle / Piece \rangle
(/Composer)
```

Example document

XQuery

XQuery is a full-fledged XML query language

```
\langle Composer \rangle
```

```
(Name) Claude Debussy (/Name)
(Vita)
  (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
   (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
  (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
   (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
(/Vita)
(Piece)
  (PTitle) La Mer (/PTitle)
  (PYear) 1905 (/PYear)
  (Instruments) Large orchestra (/Instruments)
  \langle Movements \rangle 3 \langle Movements \rangle
                                                      Example query
\langle / Piece \rangle
```

return \$x/Name

Example document

```
\langle / Composer \rangle
```

...

Schwentick

XML: Algorithms & Complexity

for \$x in doc('composers.xml')/Composer

where x/Vita/Died/Where = 'Paris'

Result

- {Name> Claude Debussy </Name>
- {Name> Eric Satie </Name>
- (Name) Hector Berlioz (/Name)
- (Name) Camille Saint-Saëns (/Name)
- {Name> Frédéric Chopin </Name>
- (Name) Maurice Ravel (/Name)
- {Name> Jim Morrison </Name>
- (Name) César Franck (/Name)
- $\langle Name \rangle$ Gabriel Fauré $\langle /Name \rangle$

(Name) George Bizet (/Name)

XQuery

XQuery is a full-fledged XML query language

ere> Paris </Where></Born>
n> Rosalie</Whom></Married>
n> Emma </Whom></Married>
> Paris </Where> </Died>

```
(Instruments) Large orcnestra (/Instruments)
(Movements) 3 (/Movements)
```

```
\langle / Piece \rangle
```

```
\langle / Composer \rangle
```

Schwentick

```
Example query
```

```
for $x in doc('composers.xml')/Composer
```

```
where x/Vita/Died/Where = 'Paris'
```

return \$x/Name

XML: Algorithms & Complexity

Transformation: XSLT

Example document

```
(Composer)
   (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     \langle Movements \rangle 3 \langle Movements \rangle
  \langle / Piece \rangle
(/Composer)
```

Transformation: XSLT

XSLT

XSLT transforms documents by means of templates

$\langle Composer \rangle$

```
(Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     \langle Movements \rangle 3 \langle Movements \rangle
  \langle / Piece \rangle
(/Composer)
```

Example documer

Transformation: XSLT

Example documer

XSLT

XSLT transforms documents by means of templates

$\langle Composer \rangle$

```
(Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
                                                          Example
     (Instruments) Large orchest
                                 (xsl:template match="Composer[Vita//Where='Paris']")
     (Movements) 3 (/Movemer
                                    (ParisComposer)
  \langle / Piece \rangle
                                       \langle xsl:copy-of select="Name"/ \rangle
                                       (xsl:copy-of select="Vita/Born"/>
(/Composer)
                                    (/ParisComposer)
                                  \langle /xsl:template \rangle
```

Transformation: XSLT __

Result	XSLT
<pre>〈ParisComposer〉</pre>	XSLT transforms documents by
(Name) Claude Debussy (/Name)	means of templates
(Born)	
〈When〉 August 22, 1862 〈/When〉	
(Where) Paris (/Where)	
⟨/Born⟩	en〉〈Where〉 Paris 〈/Where〉〈/Born〉
<pre></pre>	\\{Whom\\ Rosalie{/Whom\\/Married\\
<pre>〈ParisComposer〉</pre>	\[\] \[
(Name) Frédéric Chopin (/Name)	Where Paris (/Where (/Died)
(Born)	
〈When〉 March 1, 1810 〈/When〉	
(Where) Želazowa (/Where)	
⟨/Born⟩	Example
<pre></pre>	
<pre>〈ParisComposer〉</pre>	ate match="Composer[Vita//Where='Paris']" >
(Name) Camille Saint-Saëns (/Name)	omposer〉
(Born)	copy-of select="Name"/ \rangle
〈When〉 October 9, 1835 〈/When〉	$copy-of select = "Vita/Born" / \rangle$
(Where) Paris (/Where)	Composer
⟨/Born⟩	
<pre></pre>	hate

Schwentick

XML: Algorithms & Complexity

A Schematic View









XML: Algorithms & Complexity



Focus of this Talk

Topics

- Expressive power of XML languages
- Complexity of algorithmic tasks related to XML processing
- Tradeoff between expressiveness and complexity

Goals of this Research

- Understand expressive power and complexity of XML languages
- Identify interesting fragments with good tradeoff

Algorithmic Tasks _____

Algorithmic Tasks



Expressive power

Question: How do we measure expressive power?

Remarks

- Classes of logical formulas are a good yardstick
- \rightarrow They provide methods to prove that a query can not be expressed

Expressive power

Question: How do we measure expressive power?

Remarks

- Classes of logical formulas are a good yardstick
- \rightarrow They provide methods to prove that a query can not be expressed

Recall Relational Databases

- Core of SQL \equiv First-order Logic
- Most frequently asked queries \equiv Conjunctive queries

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Background: Complexity Classes _

Overview of Complexity Classes



Schwentick

XML, Trees and Automata _____

Question: Why is XML appealing for Theory people?



Question: Why is XML appealing for Theory people? Years ago... Theoretical Computer Science for Database Theorists: Logics, Complexity, Algorithms,... Database Theory for Theoretical Computer Scientists: terra incognita

Question: Why is XML appealing for Theory people?

Years ago...

- Theoretical Computer Science for Database Theorists: Logics, Complexity, Algorithms,...
- Database Theory for Theoretical Computer Scientists: terra incognita

After the advent of XML

Many connections between

Formal Languages & Automata Theory

and

XML & Database Theory

XML, Trees and Automata _____



Question: Why trees?

A Natural Answer

- Trees reflect the hierarchical structure of XML
- Underlying data model of XML is tree based

Question: Why trees?

A Natural Answer

- Trees reflect the hierarchical structure of XML
- Underlying data model of XML is tree based

Limitations

- But trees can not model all aspects of XML (e.g., IDREFs, data values)
- \Rightarrow Sometimes extensions are needed
 - E.g., directed graphs instead of trees

Question: Why trees?

A Natural Answer

- Trees reflect the hierarchical structure of XML
- Underlying data model of XML is tree based

Limitations

- But trees can not model all aspects of XML (e.g., IDREFs, data values)
- \Rightarrow Sometimes extensions are needed
 - E.g., directed graphs instead of trees


Question: Why trees?

A Natural Answer

- Trees reflect the hierarchical structure of XML
- Underlying data model of XML is tree based

Limitations

- But trees can not model all aspects of XML (e.g., IDREFs, data values)
- \Rightarrow Sometimes extensions are needed
 - E.g., directed graphs instead of trees



Question: Why trees?

A Natural Answer

- Trees reflect the hierarchical structure of XML
- Underlying data model of XML is tree based

Limitations

- But trees can not model all aspects of XML (e.g., IDREFs, data values)
- \Rightarrow Sometimes extensions are needed
 - E.g., directed graphs instead of trees



Nevertheless

In this tutorial we will concentrate on the

tree view at XML

Question: Why automata?

Ingredients of XML

Concepts from formal languages are obviously present around XML:

- Labelled trees
- DTD: context-free grammars
- DTD: regular expressions
- XPath: regular path expressions

We will see

- a means to define robust classes with clear semantics
- a tool for proofs
- an algorithmic tool for static analysis
- a tool for query evaluation

Question: Why automata?

Ingredients of XML

Concepts from formal languages are obviously present around XML:

- Labelled trees
- DTD: context-free grammars
- DTD: regular expressions
- XPath: regular path expressions

We will see

- \rightarrow a means to define robust classes with clear semantics
 - a tool for proofs
 - an algorithmic tool for static analysis
 - a tool for query evaluation

Question: Why automata?

Ingredients of XML

Concepts from formal languages are obviously present around XML:

- Labelled trees
- DTD: context-free grammars
- DTD: regular expressions
- XPath: regular path expressions

We will see

- a means to define robust classes with clear semantics
- \rightarrow a tool for proofs
 - an algorithmic tool for static analysis
 - a tool for query evaluation

Question: Why automata?

Ingredients of XML

Concepts from formal languages are obviously present around XML:

- Labelled trees
- DTD: context-free grammars
- DTD: regular expressions
- XPath: regular path expressions

We will see

- a means to define robust classes with clear semantics
- a tool for proofs
- \rightarrow an algorithmic tool for static analysis
 - a tool for query evaluation

Question: Why automata?

Ingredients of XML

Concepts from formal languages are obviously present around XML:

- Labelled trees
- DTD: context-free grammars
- DTD: regular expressions
- XPath: regular path expressions

We will see

- a means to define robust classes with clear semantics
- a tool for proofs
- an algorithmic tool for static analysis
- \rightarrow a tool for query evaluation

Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

From Strings to Trees _________A String

abcab

From Strings to Trees _

A String abcab





From Strings to Trees _

A String abcab



From Strings to Trees (cont.) _

XML and Trees

- XML trees are unranked : the number of children of a node is not restricted
- Automata have first been considered on ranked trees, where each symbol has a fixed number of children (rank)
- Most important ideas were already developed for ranked trees
- \rightarrow Let us take a look at this first

Remark

Sometimes trees are viewed as terms









































































Idea Tree-structured Boolean circuits Two states: q_0, q_1 Accepting at the root: q_1

 Transitions

 $\delta(\wedge, q_1) = \{(q_1, q_1)\}$
 $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$
 $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$
 $\delta(\vee, q_0) = \{(q_0, q_0)\}$
 $\delta(0, q_0) = \{\epsilon\}; \delta(0, q_1) = \emptyset$
 $\delta(1, q_1) = \{\epsilon\}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing

XML: Algorithms & Complexity



Idea Tree-structured Boolean circuits Two states: q_0, q_1 Accepting at the root: q_1

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{\epsilon\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{\epsilon\}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing



IdeaTree-structured BooleancircuitsTwo states: q_0, q_1 Accepting at the root: q_1

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{\epsilon\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{\epsilon\}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing



IdeaTree-structured BooleancircuitsTwo states: q_0, q_1 Accepting at the root: q_1

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{\epsilon\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{\epsilon\}; \delta(1, q_0) = \emptyset$



IdeaTree-structured BooleancircuitsTwo states: q_0, q_1 Accepting at the root: q_1

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{\epsilon\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{\epsilon\}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing

XML: Algorithms & Complexity


Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$

 $\delta(\land, q_1) = \{(q_1, q_1)\}$ $\delta(\land, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\lor, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\lor, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{acc\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{acc\}; \delta(1, q_0) = \emptyset$



Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$

$$\begin{split} \delta(\wedge, q_0) &= \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\} \\ \delta(\vee, q_1) &= \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\} \\ \delta(\vee, q_0) &= \{(q_0, q_0)\} \\ \delta(0, q_0) &= \{acc\}; \delta(0, q_1) = \emptyset \\ \delta(1, q_1) &= \{acc\}; \delta(1, q_0) = \emptyset \end{split}$$



Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$

 $\delta(\lor, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{ \mathsf{acc} \}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{ \mathsf{acc} \}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing



Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{acc\}; \delta(0, q_1) = \emptyset$

$$\delta(1,q_1) = \{ {\sf acc} \}; \delta(1,q_0) = \emptyset$$



Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions

$$\begin{split} \delta(\wedge, q_1) &= \{(q_1, q_1)\}\\ \delta(\wedge, q_0) &= \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}\\ \delta(\vee, q_1) &= \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}\\ \delta(\vee, q_0) &= \{(q_0, q_0)\}\\ \delta(0, q_0) &= \{acc\}; \delta(0, q_1) = \emptyset\\ \delta(1, q_1) &= \{acc\}; \delta(1, q_0) = \emptyset \end{split}$$



Idea

Guess the correct values starting from the root

Check at the leaves

Three states: $q_0, q_1,$ acc

Initial state q_1 at the root

Accepting if all leaves end in acc

Schwentick

XML: Algorithms & Complexity

Transitions $\delta(\wedge, q_1) = \{(q_1, q_1)\}$ $\delta(\wedge, q_0) = \{(q_0, q_1), (q_1, q_0), (q_0, q_0)\}$ $\delta(\vee, q_1) = \{(q_0, q_1), (q_1, q_0), (q_1, q_1)\}$ $\delta(\vee, q_0) = \{(q_0, q_0)\}$ $\delta(0, q_0) = \{acc\}; \delta(0, q_1) = \emptyset$ $\delta(1, q_1) = \{acc\}; \delta(1, q_0) = \emptyset$

Introduction - XML Processing

Regular Tree Languages

Definition

A bottom-up automaton is deterministic if

for each
$$a$$
 and $p
eq q$: $\delta(a,p) \cap \delta(a,q) = \emptyset$

Theorem

The following are equivalent for a tree language L:

- (a) L is accepted by a nondeterministic bottom-up automaton
- (b) L is accepted by a deterministic bottom-up automaton

(c) L is accepted by a nondeterministic top-down automaton

Proof idea

(a) \implies (b): Powerset construction

(a) \iff (c): Just the same thing, viewed in two different ways

Automata as Tiling Systems

Observation

• $(q_0, q_1) \in \delta(\lor, q_1)$ can be interpreted as an allowed pattern:



- A tree is accepted, iff there is a labelling with states such that
 - all local patterns are allowed
 - the root is labelled with q_1

Automata as Tiling Systems

Observation

• $(q_0, q_1) \in \delta(\lor, q_1)$ can be interpreted as an allowed pattern:



- all local patterns are allowed
- the root is labelled with q_1



 $, q_1$

Automata as Tiling Systems

Observation

• $(q_0, q_1) \in \delta(\lor, q_1)$ can be interpreted as an allowed pattern:



- A tree is accepted, iff there is a labelling with states such that
 - all local patterns are allowed
 - the root is labelled with q_1



Regular tree languages and logic

Definition: (MSO logic)
Formulas talk about
– edges of the tree (E)
– node labels $(oldsymbol{Q}_0,oldsymbol{Q}_1,oldsymbol{Q}_\wedge,oldsymbol{Q}_ee)$
 the root of the tree (root)
First-order-variables represent nodes
Monadic second-order (MSO) variables represent sets of nodes
Example: Boolean Circuits
Boolean circuit true $\equiv \exists X \; X(root) \land \forall x$
$(Q_0(x) ightarrow eg X(x)) \wedge$
$((Q_\wedge(x)\wedge X(x)) ightarrow (orall x(x,y) ightarrow X(y)]))\wedge$
$((Q_ee(x)\wedge X(x)) ightarrow (\exists y[E(x,y)\wedge X(y)]))$

Theorem [Doner 70; Thatcher, Wright 68]

 $MSO \equiv Regular Tree Languages$

Regular tree languages and logic (cont.) _

 $\frac{\text{Theorem}}{\text{MSO}} \equiv \text{Regular Tree Languages}$

Proof idea

Automata \Rightarrow MSO:

Formula expresses that there exists a correct tiling

$MSO \Rightarrow Automata$: more involved

Basic idea:

Automaton computes for each node v the set of formulas which hold

in the subtree rooted at v

Regular tree languages and logic (cont.)

Formula \Rightarrow automaton

- Let arphi be an MSO-formula, k:= quantifier-depth of arphi
- k-type of a tree t := (essentially) set of MSO-formulas ψ of quantifier-depth $\leq k$ which hold in t
- $t_1 \equiv_k t_2$: k-type $(t_1) = k$ -type (t_2)
- Automaton computes k-type of tree and concludes whether arphi holds



Det. Top-Down Automata



Det. Top-Down Automata: Acceptance

Question

What is a good acceptance mechanism for deterministic top-down automata?

Several possibilitites

(1) At all leaves states have to be accepting

(2) There is a leave with an accepting state

(2) is problematic for complement and intersection

(1) is problematic for complement and union

Det. Top-Down Automata: Acceptance (cont.)

Definition: (Root-to-frontier automata with regular acceptance condition)

- Tree automata ${\cal A}$ are equipped with an additional regular string language L over $Q imes \Sigma$
- *A* accepts *t* if the (state,symbol)-string at the leaves (from left to right) is in *L* [Jurvanen, Potthoff, Thomas 93]



A robust class

- The resulting class is closed under Boolean operations
- Good algorithmic properties
- Does not capture all regular tree languages

Schwentick

Summary

Regular tree languages

- Regular tree languages are a robust class
- Characterized by
 - parallel tree automata
 - MSO logic
 - several other models
- They are the natural analog of regular string languages
- Deterministic top-down automata with regular acceptance conditions define a weaker but nevertheless robust class

Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Tree-Walk Automata

Definition: (Tree-walk automata)

Depending on

- current state
- symbol of current node
- position of current node wrt its siblings

the automaton moves to a neighbor and takes a new state



Question What is the expressive power of tree-walk automata?

- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states




- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





- Tree-walk automata can evaluate Boolean circuit trees
- 5 states





_ A Recent Result and an Even More Recent Result __

Theorem [Bojanczyk, Colcombet 04]

Deterministic TWAs are weaker than nondeterministic TWAs

Corollary

Deterministic TWAs do not capture all

regular tree languages

Theorem [Bojanczyk, Colcombet 04]

Nondeterministic TWAs do not capture all regular tree languages

Overview of Models _



Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Decision Problems

Algorithmic problems

• We consider the following algorithmic problems

• All of them will be useful in the XML context

Membership test for ${\cal A}$	Membership test (combined)		
Given:Tree t Question:Is $t \in L(\mathcal{A})$?	Given:Automaton \mathcal{A} , tree t Question:Is $t \in L(\mathcal{A})$?		
Non-emptiness			
Given:Automaton \mathcal{A} Question:Is $L(\mathcal{A}) \neq \emptyset$?			
Containment	Equivalence		
Given:Automata $\mathcal{A}_1, \mathcal{A}_2$ Question:Is $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$?	Given: Automata $\mathcal{A}_1, \mathcal{A}_2$ Question: Is $L(\mathcal{A}_1) = L(\mathcal{A}_2)$?		

Facts

Time Bounds for the combined complexity of membership test for tree automata:

- Deterministic (parallel) tree automata: $O(|\mathcal{A}||t|)$
- Nondeterministic (parallel) tree automata: $O(|\mathcal{A}|^2|t|)$ (Compute, for each node, the set of reachable states)
- Deterministic TWAs: O(|A|²|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior function)
- Nondeterministic TWAs: O(|A|³|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior relation)

Membership Test

Facts

Time Bounds for the combined complexity of membership test for tree automata:

- \rightarrow Deterministic (parallel) tree automata: $O(|\mathcal{A}||t|)$
 - Nondeterministic (parallel) tree automata: $O(|\mathcal{A}|^2|t|)$ (Compute, for each node, the set of reachable states)
 - Deterministic TWAs: O(|A|²|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior function)
 - Nondeterministic TWAs: O(|A|³|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior relation)

Facts

Time Bounds for the combined complexity of membership test for tree automata:

- Deterministic (parallel) tree automata: $O(|\mathcal{A}||t|)$
- \rightarrow Nondeterministic (parallel) tree automata: $O(|\mathcal{A}|^2|t|)$ (Compute, for each node, the set of reachable states)
 - Deterministic TWAs: O(|A|²|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior function)
 - Nondeterministic TWAs: O(|A|³|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior relation)

Membership Test

Facts

Time Bounds for the combined complexity automata:

- Deterministic (parallel) tree automata
- Nondeterministic (parallel) tree autom (Compute, for each node, the set of reach node)
- → Deterministic TWAs: O(|A|²|t|)
 (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior function)
 - Nondeterministic TWAs: O(|A|³|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior relation)

Behavior Function

Facts

Time Bounds for the combined complexity of membership test for tree automata:

- Deterministic (parallel) tree automata: $O(|\mathcal{A}||t|)$
- Nondeterministic (parallel) tree automata: $O(|\mathcal{A}|^2|t|)$ (Compute, for each node, the set of reachable states)
- Deterministic TWAs: O(|A|²|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior function)

→ Nondeterministic TWAs: O(|A|³|t|) (Compute, for each node v, the aggregated behavior of A on its subtree: Behavior relation)

Membership Test (cont.) _

Question: What is the structural complexity for the various models?		
[Lohrey 01, Segoufin 03]		
Model	Time Complexity	Structural Complexity
Det. top-down TA	$O(\mathcal{A} t)$	LOGSPACE
Det. bottom-up TA	$O(\mathcal{A} t)$	LOGDCFL
Nondet. bottom-up TA	$O(\mathcal{A} ^2 t)$	LOGCFL
Nondet. top-down TA	$O(\mathcal{A} ^2 t)$	LOGCFL

 $O(|\mathcal{A}|^2|t|)$

 $O(|\mathcal{A}|^3|t|)$

Det. TWA

Nondet. TWA

LOGSPACE

NLOGSPACE

Facts

 Non-emptiness for string automata corresponds to Graph Reachability (complete for NLOGSPACE)

Non-emptiness for tree automata



Facts

 Non-emptiness for string automata corresponds to Graph Reachability (complete for NLOGSPACE)

Non-emptiness for tree automata



Facts

 Non-emptiness for string automata corresponds to Graph Reachability (complete for NLOGSPACE)

Non-emptiness for tree automata



Facts

 Non-emptiness for string automata corresponds to Graph Reachability (complete for NLOGSPACE)

Non-emptiness for tree automata



Facts

 Non-emptiness for string automata corresponds to Graph Reachability (complete for NLOGSPACE)

Non-emptiness for tree automata

corresponds to Path Systems



Result

- Non-emptiness for bottom-up tree
 - automata can be checked in linear time
- It is complete for **PTIME**

Containment/Equivalence _

Observations

• Of course:

 $L(\mathcal{A}_1) = L(\mathcal{A}_2) \iff [L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2) \text{ and } L(\mathcal{A}_2) \subseteq L(\mathcal{A}_1)]$

- Complexity of containment problem is very different for deterministic and non-deterministic automata
- Deterministic automata: construct product automaton

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





0110100

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





0110100

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





0|110100

Schwentick

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





01 10100

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





011 0100

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





0110 100

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





01101 00

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





011010|0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





0110100

Schwentick

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





а

b

 $\mathbf{0}$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"







Schwentick
Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





¥

а

b

 $\mathbf{0}$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





Schwentick

¥

а

b

 $\mathbf{0}$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





а

b

 $\mathbf{0}$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





Schwentick

а

 $\left(\right)$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"





Schwentick

а

 $\left(\right)$

0

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"







Schwentick

Product of 2 string automata

- "even number of zeros"
- "contains substring 00"







Schwentick

Containment: Complexity

Deterministic bottom-up tree automata

- Product automaton analogous as for string automata
 - Set of states: $Q_1 imes Q_2$
 - Transitions component-wise
- To check $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$:
 - Compute $\mathcal{B} = \mathcal{A}_1 imes \mathcal{A}_2$
 - Accepting states: $F_1 imes (Q_2-F_2)$
 - Check whether $L(\mathcal{B}) = \emptyset$
 - If so, $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$ holds

Theorem

Complexity of Containment for deterministic

bottom-up tree automata:

 $O(|\mathcal{A}_1| imes|\mathcal{A}_2|)$

Schwentick

XML: Algorithms & Complexity

Containment: Complexity (cont.)

Non-deterministic automata

- Naive approach:
 - Make \mathcal{A}_2 deterministic (size: $O(2^{|\mathcal{A}_2|}))$
 - Construct product automaton
 - \Rightarrow Exponential time

Containment: Complexity (cont.)

Non-deterministic automata

- Naive approach:
 - Make \mathcal{A}_2 deterministic (size: $O(2^{|\mathcal{A}_2|}))$
 - Construct product automaton
 - \Rightarrow Exponential time

Unfortunately...

There is essentially no better way

Containment: Complexity (cont.)

Non-deterministic automata

- Naive approach:
 - Make \mathcal{A}_2 deterministic (size: $O(2^{|\mathcal{A}_2|}))$
 - Construct product automaton

 \Rightarrow Exponential time

Unfortunately...

There is essentially no better way

Theorem [Seidl 1990]

Containment for non-deterministic tree automata is complete for **EXPTIME**

Det. Top-Down Automata: Non-Emptiness

Theorem

Nonemptiness for deterministic top-down automata \mathcal{A} can be checked in polynomial time

Proof idea

Check for each state p of \mathcal{A} and each pair (q, q') of the leaves automaton \mathcal{B} :

Is there a tree t such that \mathcal{A} starts from state p and

obtains a leave string which brings ${\cal B}$ from q to q'?



Schwentick

XML: Algorithms & Complexity

Introduction - XML Processing

Det. Top-Down Automata: Containment

Theorem

Containment for deterministic top-down automata \mathcal{A} can be checked in polynomial time

Proof idea

- Tree automata \mathcal{A}_1 , \mathcal{A}_2 with leaves automata $\mathcal{B}_1, \mathcal{B}_2$
- Check
 - for each pair (p_1,p_2) of states of \mathcal{A}_1 and \mathcal{A}_2 and
 - for each two pairs (q_1, q_1') and (q_2, q_2') of \mathcal{B}_1 and \mathcal{B}_2 , resp.:

Is there a tree t such that for both i = 1, i = 2: T_i starts from state p_i and obtains a leave string which brings \mathcal{B}_i from q_i to q'_i ?

Summary _____

Complexities of basic algorithmic problems			
Model	Membership	Non-emptiness	Containment
Det. top-down TA	LOGSPACE	ΡΤΙΜΕ	ΡΤΙΜΕ
Det. bottom-up TA	LOGDCFL	ΡΤΙΜΕ	ΡΤΙΜΕ
Nondet. bottom-up TA	LOGCFL	ΡΤΙΜΕ	EXPTIME
Nondet. top-down TA	LOGCFL	PTIME	EXPTIME
Det. TWA	LOGSPACE	PTIME (*)	PTIME (*)
Nondet. TWA	NLOGSPACE	PTIME (*)	EXPTIME (*)
(*: upper bounds)			

Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

From Ranked to Unranked Trees _



From Ranked to Unranked Trees

Agenda

• Now we move from ranked to unranked trees

- There is a basic choice:
 - Either: we encode unranked trees as binary trees and go on with ranked automata
 - Or: we adapt the ranked automata models
- In both cases: not many surprises, most results remain

Encoding Unranked Trees as Binary Trees _





Schwentick

XML: Algorithms & Complexity

Encoding Unranked Trees as Binary Trees _







Schwentick

XML: Algorithms & Complexity

_ Encoding Unranked Trees as Binary Trees (cont.) __



 \boldsymbol{e}

Binary vs. unranked trees _

Remark

- There are still other ways to encode unranked trees as binary trees
- → e.g., [Carme, Niehren, Tommasi 04]
 - We consider automata for unranked trees next

Unranked Trees: Formal Definition



XML tags can be captured by the set Σ of labels. But what about text?

- This depends on the context
- E.g., for type checking, text is irrelevant.
- In many applications, the relevant information about text nodes can be represented by predicates, e.g., whether the name = 'Debussy'.

Schwentick



Ranked trees

Transitions are described by finite sets: $\delta(\sigma, q) = \{(q_1, q_2), (q_3, q_4), \ldots\}$





Ranked trees

Transitions are described by finite sets: $\delta(\sigma, q) = \{(q_1, q_2), (q_3, q_4), \ldots\}$





Ranked treesTransitions are described by finite sets: $\delta(\sigma,q) = \{(q_1,q_2), (q_3,q_4), \ldots\}$ σ_1 σ_2 q_1 q_2



$\delta(\sigma,q)$

- For unranked trees, $\delta(\sigma,q)$ is a regular language
- δ(σ, q) can be specified by regular expression or finite string automaton
 [Brüggemann-Klein, Murata, Wood 2001]

Representation of $\delta(\sigma, q)$

Remark

- Representation of $\delta(\sigma, q)$ has influence on complexity
- Natural choice:
 - For nondeterministic tree automata: represent $\delta(\sigma, q)$ by NFAs or regular expressions
 - For deterministic tree automata: represent $\delta(\sigma, q)$ by DFAs
- \Rightarrow Same complexity results as for ranked trees

Regular sets of unranked trees

Theorem

The following are equivalent for a set L of unranked trees:

- (a) L is accepted by a nondeterministic bottom-up automaton
- (b) L is accepted by a deterministic bottom-up automaton
- (c) L is accepted by a nondeterministic top-down automaton
- (d) L is characterized by an MSO-formula

Deterministic Top-Down Automata



Checking Existence of Paths

Fact

A simple deterministic top-down automaton can check the

existence of vertical paths with regular properties

Construction

- For a node v let s(v) denote the sequence of labels from the root to v
- Let *A* be a deterministic string automaton
- $\mathcal{A}' :=$ top-down automaton which takes at v state of \mathcal{A} after reading s(v)
- \Rightarrow \mathcal{A}' is deterministic
 - There is a path from the root to a leaf vwith $s(v) \in L(\mathcal{A})$ iff \mathcal{A}' assumes at least one state from F at a leave

Streaming XML

Similar construction used for XPath evaluation on streams [Green et al. 2003]



Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Generalization of Tree-Walk Automata

Allowed transitions: Go up

Go to first child

Go to left sibling

Go to right sibling

→ Caterpillar automata [Brüggemann-Klein, Wood 2000]

Basic design choice

Should a transition to a sibling be aware of the label of the parent?



Contents

Introduction

Background on Tree Automata and Logic

Parallel Ranked Tree Automata

Sequential Ranked Tree Automata

Decision Problems for Ranked Tree Automata

Parallel Unranked Tree Automata

Sequential Unranked Tree Automata

Sequential Document Automata

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Document Automata

A third kind of automata for XML

Document automata are string automata reading XML documents as text

- Tags are represented by symbols from a given alphabet
- Variants:
 - Finite document automata
 - Pushdown document automata
- Useful especially in the context of streaming XML

Theorem [Segoufin, Vianu 02]

- Regular languages of XML-trees can be recognized by deterministic push-down document automata.
- Depth of push-down is bounded by depth of tree

Summary: Unranked Tree Automata

Summary

- Moving from ranked to unranked automata requires some adaptations
- Transitions can be defined with regular string languages $\delta(\sigma,q)$
- By and large, things work smoothly
- In particular:
 - there is an equally robust notion of regular tree languages
 - The complexities are the same as for ranked automata (if the sets $\delta(\sigma, q)$ are represented in a sensible way)

Refined Overview of Models


Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

DTDs

Specialized DTDs

1-pass Preorder Typing

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

DTDs



Some Facts

- DTDs \equiv generalized context-free grammars
- → [Berstel,Boasson 00] provide characterizations
 - Additional restriction: one-unambiguous

One-unambiguous Regular Expressions

Definition: One-unambiguous Regular Expression

- Let r be a regular expression
- $r \mapsto r'$: number the symbols of r from left to right
- $w \in L(r) \longleftrightarrow$ there is a numbered string $w' \in L(r')$

• *r* is one-unambiguous if

 $ux_iv\in L(r')$, $uy_jw\in L(r')$, $i
eq j \Rightarrow x
eq y$

Example

- $(a+b)^*ac+c \mapsto (a_1+b_2)^*a_3c_4+c_5$
- $babbac \in L(r)$ and $b_2a_1b_2b_2a_3c_4 \in L(r')$
- $(a+b)^*ac+c$ is not one-unambiguous because $b_2b_2a_3c_4\in L(r')$ and $b_2b_2a_1a_3c_4\in L(r')$
- $(b^*a)^*c$ is one-unambiguous

One-unambiguous Regular Expressions

Definition: One-unambiguous Regular Expression

- Let r be a regular expression
- $r \mapsto r'$: number the symbols of r from left to right
- $w \in L(r) \longleftrightarrow$ there is a numbered string $w' \in L(r')$

• *r* is one-unambiguous if

 $ux_iv\in L(r')$, $uy_jw\in L(r')$, $i
eq j \Rightarrow x
eq y$

Example

- $(a+b)^*ac+c \mapsto (a_1+b_2)^*a_3c_4+c_5$
- $babbac \in L(r)$ and $b_2a_1b_2b_2a_3c_4 \in L(r')$
- $(a+b)^*ac+c$ is not one-unambiguous because $b_2b_2 \ a_3 \ c_4 \in L(r')$ and $b_2b_2 \ a_1 \ a_3c_4 \in L(r')$
- $(b^*a)^*c$ is one-unambiguous

One-unambiguous Regular Expressions

Definition: One-unambiguous Regular Expression

- Let r be a regular expression
- $r \mapsto r'$: number the symbols of r from left to right
- $w \in L(r) \longleftrightarrow$ there is a numbered string $w' \in L(r')$

• *r* is one-unambiguous if

 $ux_iv\in L(r')$, $uy_jw\in L(r')$, $i
eq j \Rightarrow x
eq y$

Example

- $(a+b)^*ac+c \mapsto (a_1+b_2)^*a_3c_4+c_5$
- $babbac \in L(r)$ and $b_2a_1b_2b_2a_3c_4 \in L(r')$
- $(a+b)^*ac+c$ is not one-unambiguous because $b_2b_2 \ a_3 \ c_4 \in L(r')$ and $b_2b_2 \ a_1 \ a_3c_4 \in L(r')$
- $(b^*a)^*c$ is one-unambiguous

Restriction

- Expressions in DTDs have to be one-unambiguous
- Inherited from SGML

Validation wrt a DTD



Validation wrt a DTD (cont.)

Observation

- Validation wrt DTDs is a very simple task
- Can be done by
 - Bottom-up automata
 - Deterministic top-down automata (if siblings contribute to new state)
 - Deterministic tree-walk automata:
 Just a depth-first left-to-right traversal
- In particular: Validation possible in linear time during one pass through the document

1-pass validation)

• Further: DTDs are always satisfiable

Containment of DTDs



 \Rightarrow Containment of regular expressions r_1, r_2 by product automaton of size $O(|r_1||r_2|)$

Containment of DTDs (cont.)

Question

What if the requirement of being one-unambiguous is dropped?

A classical result

Theorem [Stockmeyer, Meyer 71]

Containment and Equivalence for regular expressions on strings are complete for **PSPACE**

Corollary

Containment of DTDs (with unrestricted regular expressions) is **PSPACE**-complete

Theorem [Martens, Neven, Sch. 04]

Containment and Equivalence for regular expressions are

- **coNP**-complete for concatenations of a, b, c and a^*, b^*, c^*
- **coNP**-complete for concatenations of *a*, *b*, *c* and *a*?, *b*?, *c*?
- **PSPACE**-complete for concatenations of a, b, c and

 $(a^* + b^* + \dots + c^*)$

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

 DTDs

Specialized DTDs

1-pass Preorder Typing

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Weakness of DTDs



Observation

- Elements with the same name may have different structure in different contexts
- \rightarrow It would be nice to have types for elements

Specialized DTDs

Schwentick

mode

ad

model

year

Specialized DTDs



Note

Concerning the name:

"specialized" refers to types, not to DTDs

Example

Dealer \rightarrow UsedCars NewCars μ (Dealer) = Dealer UsedCars \rightarrow adUsed* NewCars \rightarrow adNew^{*} adUsed \rightarrow model year $adNew \rightarrow model$

 μ (UsedCars) = UsedCars μ (NewCars) = NewCars μ (adUsed) = ad μ (adNew) = ad

Schwentick

A Further Example ____

Example: SDTD for Boolean circuit trees					
	$(1 \cap \mathbf{P} \mid 1 \wedge \mathbf{N} \mathbf{D} \mid 1 \mid \mathbf{h} \in \mathbf{f})$	Tag	μ (Tag)		
I-AND -	$\rightarrow (1 - OK 1 - AND 1 - Ieal)^{-1}$	1-AND	AND		
I-OR –	\rightarrow .* (I-OR I-AND I-leaf) .*	0-AND	AND		
0-AND -	\rightarrow .* (0-OR 0-AND 0-leaf) .*	1-OR	OR		
0-OR -	\rightarrow (0-OR 0-AND 0-leaf)*	0-OR	OR		
1-leaf -	$ ightarrow \epsilon$		1		
0-leaf –	$\rightarrow \epsilon$	1-leal	L		
		0-leaf	0		

Specialized DTDs (cont.) _

Observation

- A naive validation by exhaustively trying all possible functions μ requires exponential time
- But help comes from automata...

Specialized DTDs (cont.)

Observation

- A naive validation by exhaustively trying all possible functions μ requires exponential time
- But help comes from automata...
- A tree conforms to a specialized DTD (d, μ) if there is a labeling of its nodes by types which is valid wrt. d

Specialized DTDs (cont.)

Observation

- A naive validation by exhaustively trying all possible functions μ requires exponential time
- But help comes from automata...
- A tree conforms to a specialized DTD (d, μ) if there is a labeling of its nodes by types which is valid wrt. d
- This reminds us of something...

Specialized DTDs (cont.)

Observation

- A naive validation by exhaustively trying all possible functions

 µ requires exponential time
- But help comes from automata...
- A tree conforms to a specialized DTD (d, μ) if there is a labeling of its nodes by types which is valid wrt. d
- This reminds us of something...

Theorem

Specialized DTDs capture exactly the regular tree languages

Validation and Typing

Definition: (Validation)				
Given: Specialized DTD d , tree t				
Qeustion: Is t valid wrt d ?				
Definition: (Typing)				
Given: Specialized DTD d , tree t				
Given. Specialized DTD u , the i				

Facts

- Specialized DTDs \equiv regular tree languages
- \rightarrow Validation in linear time by deterministic push-down automata
 - Typing in linear time (Bottom-up automaton) the document
 - Satisfiability \equiv Non-emptiness of tree automata: **PTIME**

Restricted Schemas

- ullet Two types b,b' compete if $\mu(b)=\mu(b')$
- A specialized DTD is single-type if no competing types occur in the same rule $(e.g., a \rightarrow bcb' \text{ is not single-type})$
- A specialized DTD is restrained-competition if no rule allows strings wbv, wb'v' with competing types b, b'(e.g., $a \rightarrow c(b + d^*b')$ is not restrained-competition)
- The authors argue that XML-Schema roughly corresponds to single-type SDTDs

Restricted Schemas

$$\succ$$
 Two types b,b' compete if $\mu(b)=\mu(b')$

- A specialized DTD is single-type if no competing types occur in the same rule $(e.g., a \rightarrow bcb' \text{ is not single-type})$
- A specialized DTD is restrained-competition if no rule allows strings wbv, wb'v' with competing types b, b'(e.g., $a \rightarrow c(b + d^*b')$ is not restrained-competition)
- The authors argue that XML-Schema roughly corresponds to single-type SDTDs

Restricted Schemas

- ullet Two types b,b' compete if $\mu(b)=\mu(b')$
- A specialized DTD is single-type if no competing types occur in the same rule $(e.g., a \rightarrow bcb' \text{ is not single-type})$
 - A specialized DTD is restrained-competition if no rule allows strings wbv, wb'v' with competing types b, b'(e.g., $a \rightarrow c(b + d^*b')$ is not restrained-competition)
 - The authors argue that XML-Schema roughly corresponds to single-type SDTDs

Restricted Schemas

- Two types b,b' compete if $\mu(b)=\mu(b')$
- A specialized DTD is single-type if no competing types occur in the same rule $(e.g., a \rightarrow bcb' \text{ is not single-type})$
- A specialized DTD is restrained-competition if no rule allows strings wbv, wb'v' with competing types b, b'(e.g., $a \rightarrow c(b + d^*b')$ is not restrained-competition)
 - The authors argue that XML-Schema roughly corresponds to single-type SDTDs

Restricted Schemas

(Murata, Lee, Mani 2001) introduced* restrictions on specialized DTDs to ensure efficient validation (*: in a slightly different framework)

- Two types b,b' compete if $\mu(b)=\mu(b')$
- A specialized DTD is single-type if no competing types occur in the same rule $(e.g., a \rightarrow bcb' \text{ is not single-type})$
- A specialized DTD is restrained-competition if no rule allows strings wbv, wb'v' with competing types b, b'

(e.g., $a \rightarrow c(b + d^*b')$ is not restrained-competition)

→ The authors argue that XML-Schema roughly corresponds to single-type SDTDs

Schema Containment

Schema Containment

Given: Schemas d_1, d_2

Question: Is $L(d_1) \subseteq L(d_2)$?

Observations

- Important, e.g., for data integration
- Recall: Specialized DTDs are essentially non-deterministic tree automata
- \Rightarrow Containment of specialized DTDs is in **EXPTIME**
 - But the restricted forms have lower complexity
 - Complexity of containment depends on the allowed regular expressions

Schema Containment: Complexity

Results (partly from [Martens, Neven, Sch. 04])				
	Schema type	unrestricted	deterministic	
			expressions	
	DTDs	PSPACE	PTIME	
	single-type SDTDs	PSPACE	ΡΤΙΜΕ	
	restrained-competition	PSPACE	ΡΤΙΜΕ	
	SDTDs			
	unrestricted SDTDs	EXPTIME	EXPTIME	

Observations

- For unrestricted SDTDs the complexity is dominated by tree automata containment
- For the others it is dominated by the sub-task of checking containment for regular expressions

Schema Containment: Complexity

Observations (cont.)

- ... for the others it is dominated by the sub-task of checking containment for regular expressions
- Actually, this observation can be made more precise

Theorem [Martens, Neven, Sch. 04]

For a class \mathcal{R} of regular expressions and a complexity class \mathcal{C} , the following are equivalent

- (a) The containment problem for \mathcal{R} expressions is in \mathcal{C} .
- (b) The containment problem for DTDs with regular expressions from \mathcal{R} is in \mathcal{C} .
- (c) The containment problem for single-type SDTDs with regular expressions from \mathcal{R} is in \mathcal{C} .

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

 DTDs

Specialized DTDs

1-pass Preorder Typing

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Typing (cont.)

Observations

- Type of a node \equiv state of deterministic bottom-up automaton
- Deterministic push-down automaton can assign types during 1 pass
- But the type of a node v is determined after visiting its subtree



Typing (cont.)

Observations

- Type of a node \equiv state of deterministic bottom-up automaton
- Deterministic push-down automaton can assign types during 1 pass
- But the type of a node v is determined after visiting its subtree
- 1-pass preorder typing

determine type of v before visiting the subtree of v





1-Pass Preorder Typing

Question

When would it be important to know the type of \boldsymbol{v} before visiting the subtree of \boldsymbol{v} ?



1-Pass Preorder Typing



Answer

Whenever the processing proceeds in document order, e.g.:

- Streaming XML: Typing as the first operator in a pipeline
- SAX-based processing

1-Pass Preorder Typing



Our next goal

Find out which schemas admit 1-pass preorder typing

1-Pass Preorder Typing (cont.)

Remarks

- The definition of "1-pass preorder typing" does not yet restrict the efficiency of determining the type of a node
- Typing could be 1-pass preorder but very time consuming
- It turns out that essentially this never happens
- Clearly, restrained competition is sufficient for 1-pass preorder typing
- Is it also necessary?

1-Pass Preorder Typing (cont.)

Remarks

- The definition of "1-pass preorder typing" does not yet restrict the efficiency of determining the type of a node
- Typing could be 1-pass preorder but very time consuming
- It turns out that essentially this never happens
- Clearly, restrained competition is sufficient for 1-pass preorder typing
- Is it also necessary?

Theorem [Martens, Neven, Sch. 2004]

For a regular tree language L the following are equivalent

- (a) L can be described by a 1-pass preorder typable SDTD
- (b) L can be described by a restrained-competition SDTD
- (c) L has linear time 1-pass pre-order typing
- (d) L can be preorder-typed by a deterministic pushdown document automaton
- (e) Types for trees in L can be computed by a left-siblings-aware top-down deterministic tree automaton

A Very Robust Class

Further characterizations

- This class has further interesting characterizations
- E.g., by closure under ancestor-sibling-guarded subtree exchange



Theorem [Martens, Neven, Sch. 2004]

For a regular tree language L the following are equivalent

- (a) L can be described by a single-type SDTD
- (b) Types for trees in L can be computed by a simple top-down deterministic tree automaton
- (c) L is closed under ancestor-guarded subtree exchange
Summary: Schema Languages

Summary

Expressive power

- Regular tree languages offer a nice framework (\equiv MSO logic!)
- Restrained competition \equiv Deterministic top-down automata

Validation Linear time

Typing

- Linear time
- Efficient 1-pass preorder typing for restrained competition SDTDs

Satisfiability

- DTDs: trivial
- SDTDs: **PTIME**

Containment

- General SDTDs: **EXPTIME**
- Restrained competition SDTDs: PTIME

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

Node-Selecting Queries



Node-Selecting Queries



Node-Selecting Queries (cont.)

Question

Is there a class of node-selecting queries, as robust

as the regular tree languages?

Node-Selecting Queries (cont.)

Question

Is there a class of node-selecting queries, as robust as the regular tree languages?

Observation

- There is a simple way to define node selecting queries by monadic second-order formulas:
- Simply use one free variable: $\varphi(x)$
- Is there a corresponding automaton model?
- It is relatively easy to add node selection to nondeterministic bottom-up automata

Node-Selecting Queries (cont.)

Question

Is there a class of node-selecting queries, as robust as the regular tree languages?

Observation

- There is a simple way to define node selecting queries by monadic second-order formulas:
- Simply use one free variable: $\varphi(x)$
- Is there a corresponding automaton model?
- It is relatively easy to add node selection to nondeterministic bottom-up automata

Definition: (Nondetermistic bottom-up node-selecting automata)

• Nondeterministic bottom-up automata plus select function:

$s:Q imes \Sigma o \{0,1\}$

• Node v is in result set for tree $t:\iff$ there is an accepting computation on t in which v gets a state q such that $s(q,\lambda(v))=1$



้อ

- $\epsilon + q_0^* + Q^* q_a Q^*$
- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

C

e



้อ

 q_b

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

 q_0

C

 q_0

e



 q_a

h

้ล

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

 q_0

 q_b

С

C

 q_b

 q_0

e

 q_0

 q_0



 q_a

h

้ล

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

 q_0

h

 q_0

e

 q_0

 q_0

 q_b

C

C

 q_b



 q_b

h

• Accepting: q_0

 q_0

e

 q_0

C



้อ

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

C

e



้อ

 q_b

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

 q_0

C

 q_0

e

 q_b

h



้ล

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

 q_0

 q_b

C

C

 q_b

h

 q_0

е

 q_0

 q_b

h

 q_b



้ล

- all other sets empty
- $s(q_b, b) = 1$
- all others: 0
- Accepting: q_0

Schwentick

 q_0

h

 q_0

e

 q_0

 q_b

h

 q_b

C

C

 q_b

h

 q_b



Node-Selecting Automata

Fact

- Existential semantics: a node is in the result if there exists an accepting run which selects it
- Universal semantics: a node is in the result if every accepting run selects it
- Both semantics define the same class of queries

Result

A node selecting query is MSO-definable iff it is expressible by a nondeterministic bottom-up node selecting automaton

Node-Selecting Automata (cont.)

Result

A node selecting query is MSO-definable iff it is expressible by a nondeterministic bottom-up node selecting automaton

Proof idea

- Given formula $\varphi(x)$ of quantifier-depth k and tree t, for each node v the automaton does the following:
 - Compute k-type of subtree at v
 - Guess k-type of "envelope tree" at v
 - Conclude whether $oldsymbol{v}$ is in the output
 - Check consistency upwards towards the root
- \Rightarrow one unique accepting run



97

Equivalent Models

More query models

- Unfortunately, the translation from formula to automaton can be prohibitively expensive: number of states $\sim 2^{2^{|\varphi|}} \left. \left. \right\} |\varphi|$
- Actually: If P ≠ NP there is no elementary *f*, such that MSO-formulas can be evaluated in time *f*(|formula| × *p*(|tree|)) with polynomial *p* [Frick, Grohe 2002]
- \rightarrow query languages with better complexity properties needed
 - Good candidate: Monadic Datalog [Gottlob, Koch 2002] and its restricted dialects like TMNF
 - Further models:
 - Attributed Grammars [Neven, Van den Bussche 1998]
 - μ -formulas [Neumann 1998]
 - Context Grammars [Neumann 1999]
 - Deterministic Node-Selecting Automata [Neven, Sch. 1999]

Node-selecting Queries: Evaluation Complexity

Some facts about query evaluation

- MSO node-selecting queries can be evaluated in two passes through the tree
 - first pass, bottom-up: essentially computes the types of the subtrees
 - second pass, top-down: essentially computes the types of the envelopes and combines it with the subtree information
- This can be implemented by a 2-pass pushdown document automaton which in its first pass attaches information to each node [Neumann, Seidl 1998; Koch 2003]
- In particular: queries can be evaluated in linear time

Node-selecting Queries: Static Analysis

Facts

- Satisfiability: Non-emptiness of node-selecting automata is **PTIME**-complete
- Satisfiability of MSO-queries is non-elementary
- Containment of node-selecting automata is
 EXPTIME-complete

Summary: Node-Selecting Queries

Summary There is a natural notion of regular node-selecting queries generalizing regular tree languages

- Probably for most practical purposes too strong
- But it offers a useful framework for the study of other classes of queries
- A robust but weaker class of queries is captured by pebble automata

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

XPath Fragments

XPath and Important Fragments

- Many fragments of XPath have been defined
- The main fragments we consider are:
 - Full XPath : XPath 1.0

(besides the namespace related functions)

 Navigational XPath [Gottlob, Koch, Pichler 03, Benedikt, Fan, Kuper 03]:

Location paths along all axes plus Boolean operations (no attributes, no relational operators)

- Forward XPath: Navigational XPath restricted to child, descendant, self, descendant-or-self

XPath

Main Ingredients of Navigational XPath

• Location Step :

 $p = Axis :: Node-Test Predicate^*$

- Predicate : [Expression]
- Location Path :

 $\pi =$ Location Step / Location Path

More explicitly: $\pi = p_1 / \cdots / p_k$

• Expression : basically a Boolean combination of location steps

Example

/descendant::a/

child::*[descendant::c and not following-sibling::b]/

descendant::a

XPath Example

Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/desc::a



XPath Example

Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/desc::a



XPath Semantics

- Result of an expression is a node set or a single value (Boolean, number or string)
- Expressions are evaluated relative to a **context**, in particular relative to a **context node**
- Location step: p = (a :: n q) relative to context node u yields the set [p](u) of nodes v such that
 - (u, v) are in a-relation
 - v is labeled according to n (arbitrary, if n = *)
 - all predicates of ${m q}$ hold at ${m v}$
- Extended to sets S of nodes: $\llbracket p \rrbracket(S) = \cup_{u \in S} \llbracket p \rrbracket(u)$
- Location path: $\llbracket p/\pi \rrbracket(S) = \llbracket \pi \rrbracket(\llbracket p \rrbracket(S))$

Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/desc::a



Example XPath Expression

/desc::a /child::*[desc::c and not foll-sib::b]/desc::a



Example XPath Expression

/desc::a/ child::* [desc::c and not foll-sib::b]/desc::a



Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/desc::a



Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/desc::a



Example XPath Expression

/desc::a/child::*[desc::c and not foll-sib::b]/ desc::a



XPath Expression as Query Tree


A simplified Notation [Benedikt, Fan, Kuper 03]

Notation

• \downarrow , \uparrow , \rightarrow , \leftarrow , \circlearrowleft : child, parent, next-sibling, previous-sibling, self

•
$$\downarrow^+$$
, \uparrow^+ , \rightarrow^+ , \leftarrow^+ :

descendant, ancestor, following-sibling, preceding-sibling

•
$$\downarrow^*$$
, \uparrow^* , \rightarrow^* , \leftarrow^* :

descendant-or-self, ancestor-or-self,
following-sibling-or-self, preceding-sibling-or-self

Example

- child::a/descendant::c/following-sibling::*/parent::b can be expressed as $\downarrow/a/\downarrow^+/c/\rightarrow/\uparrow/b$
- The following-axis can be expressed via $\uparrow^*/\rightarrow^+/\downarrow^*$

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

Xpath and Existential First-order Logic

Characterizations of XPath [Benedikt, Fan, Kuper 03]

- Navigational XPath (without not and and) corresponds to positive existential first-order logic
- Different XPath axes correspond to different signatures

Proof idea

• Basic idea:

For each node v of the query tree: guess a node h(u) in the document tree and check that h is a "homomorphism"

• Main difficulty in proof:

Deal with conjunctions of conditions

Further Results on

- closure properties
- axiomatizations of equivalence

Backward vs. Forward Axes

Elimination of Backward Axes [Olteanu et al. 02]

- In absolute XPath expressions, all backward axes can be eliminated
- Two sets of rewrite rules:
 - introducing equality expressions, linear time (and size)
 - without equality expressions, possibly exponential size

Xpath and First-Order Logic

Reminder

Navigational XPath without negation corresponds to positive existential first-order logic

Question: What is needed to capture full first-order logic?

Xpath and First-Order Logic

Reminder

Navigational XPath without negation corresponds to positive existential first-order logic

Question: What is needed to capture full first-order logic?

Conditional axes

Conditional axes

Expressions of the kind P^+ , where P is an expression

Example

 $(\texttt{child} :: a[\texttt{desc} :: b \text{ or } \texttt{child} :: c])^+$

holds between $oldsymbol{u}$ and $oldsymbol{v}$ if

- ullet v is a descendant of u and
- all intermediate nodes
 - are labelled with $oldsymbol{a}$ and
- have a *c*-child or a *b*-descendant XML: Algorithms & Complexity Introduction - XPath

Xpath and First-Order Logic (cont.)

Theorem [Marx 04]

Navigational XPath with conditional axes corresponds exactly to first-order logic (wrt node-selecting queries)

Proof idea

The proof uses a decomposition technique similar to the proof that LTL corresponds to first-order logic over linear structures [Gabbay et al. 80]

Pebble Automata and XPath

Definition: Pebble Automata

- Extension of tree-walk automata by fixed number of pebbles
- Only pebble with highest number (current pebble) can move, depending on state, number of pebbles symbols under pebbles and incidence of pebbles
- Possible pebble movements:
 - stay, go to left sibling, go to right sibling, go to parent
 - lift current pebble or place new pebble at current position
- Nondeterminism possible

Fact Deterministic pebble automata capture navigational XPath queries Proof idea For each node of the query tree: cycle through all possible nodes of the document tree

Automata and Logic

Some observations

- On strings, MSO logic and (unary) transitive closure logic (TC-logic) coincide
- On trees
 - MSO \equiv parallel automata
 - TC-logic \equiv pebble automata (i.e., strongest sequential automata)
- Whether on trees MSO \equiv TC-logic is open
- The relationship between logics and automata models between FO and TC-logic is largely unexplored:
 - Tree-walk automata,
 - FO-logic + regular expressions
 - Conditional XPath + arbitrary star operator

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

XPath Query Evaluation

Naive Evaluation

```
Procedure \mathsf{Eval}(p_1/\cdots/p_n,v)
```

```
S:=[\![p_1]\!]v
```

```
\mathsf{IF}\; n = 1 \; \mathsf{RETURN}\; S\; \mathsf{ELSE}\; S' := \emptyset
```

```
FOR u \in S DO S' := S' \cup \mathsf{Eval}(p_2 / \cdots / p_n, u)
```

RETURN S'

XPath Query Evaluation

Naive Evaluation

Procedure $\operatorname{Eval}(p_1/\cdots/p_n,v)$

 $S := \llbracket p_1
rbracket v$

 $\mathsf{IF} \ n = 1 \ \mathsf{RETURN} \ S \ \mathsf{ELSE} \ S' := \emptyset$

FOR $u \in S$ DO $S' := S' \cup \mathsf{Eval}(p_2 / \cdots / p_n, u)$

RETURN S'

Complexity

- $T(p_1/\cdots/p_n,t) = O(\text{size of } t) \times T(p_2/\cdots/p_n,t)$
- Could be exponential
- Experiments (reported in [Gottlob, Koch, Pichler 02]) show that available XPath processors had exponential complexity



Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]



Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]





Basic Idea

Combine top-down evaluation of the "main path" with bottom-up evaluation of predicates [Gottlob, Koch, Pichler 02]









Complexity

- For each node of the query tree: O(|t|) steps
- Overall: O(query size $\times |t|)$

Beyond Navigational XPath

Example expression

/desc::a/child::*[desc::c[position() > 1]]/desc::a

Observations

- In general, a subexpression does not only depend on a context node but also on
 - context position (position())
 - context size (last())
- \rightarrow predicates can no longer be evaluated in a bottom-up fashion
 - Basic idea of [Gottlob, Koch, Pichler 02]: Compute the value of each subexpression for each triple (v, i, l) of
 - a node v
 - a position i

– a size *l*

Two Algorithms for XPath Evaluation

Results from [Gottlob, Koch, Pichler 02/03]

- The basic idea can be turned into different algorithms:
 - a bottom-up algorithm:
 - * Computing the value for each e, (v, i, l) in a dynamic programming fashion
 - * Time bound: $O((\text{tree size})^5 \times (\text{query size})^2)$
 - a (mixed) top-down algorithm:
 - * Compute as much information as possible in top-down fashion to evaluate subexpressions only for relevant triples (v, i, l)
 - * Time bound: $O((\text{tree size})^4 \times (\text{query size})^2)$
- Further time bound for the "extended Wadler fragment": $O((\text{tree size})^2 \times (\text{query size})^2)$

Structural Complexity of XPath Evaluation

Further Results

- In [Gottlob, Koch, Pichler 03] the complexity of XPath evaluation is considered
- Data Complexity:
 - Navigational XPath: LOGSPACE-complete (e.g., via pebble automata)
 - Full XPath: also LOGSPACE (?)
- Combined Complexity:
 - Navigational XPath: **PTIME**-complete
 - Positive Navigational XPath: LOGCFL-complete
 - An even much larger fragment (pXPath) is in LOGCFL

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

XPath satisfiability

Observation

Not all XPath expressions are satisfiable, e.g.:

child::a/child::b/following-sibling::c/parent::d

Question

What is the complexity of checking satisfiability of an XPath

expression for different fragments?

XPath satisfiability

Observation

Not all XPath expressions are satisfiable, e.g.:

child::a/child::b/following-sibling::c/parent::d

Question

What is the complexity of checking satisfiability of an XPath expression for different fragments?

Theorem [Hidders 03]

- Satisfiability for positive navigational XPath expressions is in NP
- Even for expressions without Boolean operators it is NP-hard
- For relative expressions without Boolean operators it is in P

XPath satisfiability

Observation

Not all XPath expressions are satisfiable, e.g.:

child::a/child::b/following-sibling::c/parent::d

Question

What is the complexity of checking satisfiability of an XPath expression for different fragments?

Theorem [Hidders 03]

- Satisfiability for positive navigational XPath expressions is in NP
- Even for expressions without Boolean operators it is NP-hard
- For relative expressions without Boolean operators it is in P

Remark

As navigational XPath can express star-free regular expressions along a path:

Satisfiability of navigational XPath is non-elementary

Theorem [Hidders 03]

Satisfiability for positive navigational XPath expressions is in NP

Proof idea

- If an expression e without ∪ is satisfiable it has a model of size ≤ |e|
- For an arbitrary (negation-free) expression guess a disjunct of the disjunctive normal form

Theorem [Hidders 03]

Satisfiability for positive navigational XPath expressions without Boolean operators is NP-hard

Proof idea

- Reduction from *Bounded Multiple String Matching (BMS)*:
 - Given: Pattern strings p_1, \ldots, p_n over $\{0, 1, *\}$
 - Question: Is there a string over $\{0,1\}$ of length $|p_1|$ which matches all patterns?
- Example: *0**1, 00*1, *111 has solution 00111
- As XPath expression:
 /↓/↓/0/↓/↓/1 [↑*/1/↑/0/↑/0] [↑*/1/↑/1/↑/1]

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

Node-selecting Queries

XPath: Semantics and Fragments

XPath: Expressive Power

XPath: Evaluation

XPath: Satisfiability

XPath: Containment

XSLT

XQuery

Conclusion

Forward XPath

Example query //Vita/Died/*

Example document

```
\langle Composer \rangle
```

```
 (Name) Claude Debussy (/Name)
(Vita)
  (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
  (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
  (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
  (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
(/Vita)
(Piece)
  (PTitle) La Mer (/PTitle)
  (PYear) 1905 (/PYear)
  (Instruments) Large orchestra (/Instruments)
  \langle Movements \rangle 3 \langle Movements \rangle
```

 $\langle / \mathsf{Piece} \rangle$

```
\langle / Composer \rangle
```

Forward XPath

Example query //Vita/Died/*

Example document

```
(Composer)
  (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     \langle Movements \rangle 3 \langle Movements \rangle
     . . .
  〈/Piece〉
(/Composer)
```

Abbreviated Syntax for Forward XPath

Example document

```
(Composer)
  (Name) Claude Debussy (/Name)
  (Vita)
     (Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)
     (Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
     (Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
     (Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
  (/Vita)
  (Piece)
     (PTitle) La Mer (/PTitle)
     (PYear) 1905 (/PYear)
     (Instruments) Large orchestra (/Instruments)
     (Movements) 3 (/Movements)
  \langle / Piece \rangle
  . . .
(/Composer)
```
Abbreviated Syntax for Forward XPath

Another example query

(/*[Name]//When) | (//Where)

Example do

```
(Composer)
```

```
(Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)
(Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)
(Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)
```

```
{Died}{When} March 25, 1918 {/When}{Where} Paris {/Where} {/Died}
{/Vita}
```

```
(Piece)
```

```
〈PTitle〉 La Mer 〈/PTitle〉
〈PYear〉 1905 〈/PYear〉
〈Instruments〉 Large orchestra 〈/Instruments〉
```

```
\langle Movements \rangle 3 \langle /Movements \rangle
```

```
\langle / \mathsf{Piece} \rangle
```

```
\langle / Composer \rangle
```

Abbreviated Syntax for Forward XPath

Another example query

(/*[Name]//When) | (//Where)

Example do

(Composer)

```
{Name> Claude Debussy </Name>

(Vita>)
```

 $\langle Born \rangle \langle When \rangle August 22, 1862 \langle When \rangle \langle Where \rangle Paris \langle Where \rangle \langle Born \rangle$

(Married) (When) October 1899 (/When) (Whom) Rosalie (/Whom) (/Married)

(Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)

(Died) (When) March 25, 1918 (/When) (Where) Paris (/Where) (/Died)

$\langle / Vita angle$

(Piece)

```
{PTitle> La Mer </PTitle>
<PYear> 1905 </PYear>
<Instruments> Large orchestra </Instruments>
<Movements> 3 </movements>
```

```
\langle / Piece \rangle
```

```
\langle / Composer \rangle
```

More XPath operators

Operator	Meaning
p/q	child
p//q	descendant
p[q]	filter
*	wildcard
$p \mid q$	disjunction

Abbreviated Syntax for Forward XPath

Another example query

(/*[Name]//When) | (//Where)

(Composer)

```
 (Name Claude Debussy (/Name)
```

(Vita)

(Born) (When) August 22, 1862 (/When) (Where) Paris (/Where) (/Born)

(Married) (When) October 1899 (/When) (Whom) Rosalie(/Whom) (/Married)

(Married) (When) January 1908 (/When) (Whom) Emma (/Whom) (/Married)

(Died) (When) March 25, 1918 (/When) (Where) Paris (More VD

Example do

$\langle /Vita \rangle$

(Piece)

```
{PTitle> La Mer </PTitle>
<PYear> 1905 </PYear>
<Instruments> Large orchestra </Instruments>
<Movements> 3 </movements>
```

```
\langle / Piece \rangle
```

```
(/Composer)
```

More XPath operators

Operator	Meaning
p/q	child
p//q	descendant
p[q]	filter
*	wildcard
$p \mid q$	disjunction









Counter-example	
<pre>〈Vita〉</pre>	
<pre>〈Died〉</pre>	
(How) Heart disease (/How)	
$\langle / Died \rangle$	





Further question

But what if the type of documents is constrained?

XPath Containment (cont.) ____

Fact
For all XML documents of type
Composers [</td
ELEMENT Composers (Composer*)
ELEMENT Composer (Name, Vita, Piece*)
ELEMENT Vita (Born, Married*, Died?)
ELEMENT Born (When, Where)
ELEMENT Married (When, Whom)
ELEMENT Died (When, Where)
ELEMENT Piece (PTitle, PYear,</td
Instruments, Movements)>
]>
the pattern <mark>//Vita/Died/*</mark> always selects a subset of positions of
<pre>(/*[Name]//When) (//Where)</pre>

XPath Containment: Definition

Definition: Containment for XPath(S)

Let S be a set of XPath-operators. The containment problem for XPath(S) is:

Given: XPath(S)-expression p, q

Question: Is $p(t) \subseteq q(t)$ for all documents t?

Definition: Containment for XPath (S) with DTD

Let S be a set of XPath-operators. The containment problem for XPath(S) in the presence of DTDs is:

Given: XPath(S)-expression p, q, DTD d

Question: Is $p(t) \subseteq q(t)$ for all documents t satisfying $t \models d$?

Observation

These problems are crucial for static analysis and query optimization

Question For which fragments *S* are these problems • decidable?

• efficiently solvable?

General remarks

- The XPath containment problem has been considered for various sets of operators
- Focus on Forward XPath
- Results vary from **PTIME** to "undecidable"
- Various methods have been used:
 - Canonical model technique
 - Homomorphism technique
 - Chase technique
- More about this in [Miklau, Suciu 2002; Deutsch, Tanen 2001; Sch. 2004]
- We will consider automata based techniques

The Automata Technique

Definition: (Relative Containment for XPath (S) wrt DTD) Let S be a set of XPath-operators. The containment problem for XPath(S) relative to a DTD is: Given: XPath(S)-expression p, q, DTD dQuestion: Is $p(D) \subseteq q(D)$ for all documents D satisfying $D \models d$?

A vague plan

- Construct an automaton \mathcal{A}_p for p
- Construct an automaton \mathcal{A}_q for q
- Construct an automaton \mathcal{A}_d for d
- Combine these automata suitably to get an automaton which accepts all counter-example documents

A Simplification

Definition: (Boolean containment)

 $p \subseteq_b q :\iff$ whenever p selects *some* node in a

tree t then q also selects some node in t.

Useful observation [Miklau, Suciu 2002]

In the presence of [], Boolean containment has the same complexity as containment.



XML: Algorithms & Complexity

Introduction - XPath

XPath Containment: 2 Examples

Result 1 [Neven, Sch. 2003]

The Boolean containment problem for XPath(/, //) in the presence of DTDs is in **PTIME**

Result 2 [Neven, Sch. 2003]

The Boolean containment problem for XPath(/, //, [], *, |) in the presence of DTDs is in **EXPTIME**

Note

Both results are optimal wrt complexity:

the problems are complete for these classes

Containment for XPath(/,//) and DTDs

Result 1 [Neven, Sch. 2003]

The Boolean containment problem for XPath(/, //) in the presence of DTDs is in **PTIME**

Proof idea

- XPath(/, //)-expressions can only describe vertical paths in a tree
- Each expression is basically of the form $p_1//p_2//\cdots//p_k$, where each p_i is of the form $l_{i1}/\cdots/l_{im_i}$
- $\bullet\,$ On strings this is a sequence of string matchings corresponding to a regular language L
- \Rightarrow Deterministic string automaton of linear size
 - Recall: there is a deterministic top-down automaton which checks whether a *p*-path exists
- \Rightarrow Deterministic top-down automaton \mathcal{A}_p
- \Rightarrow Deterministic top-down automaton $\mathcal{A}_{\overline{q}}$ checking that no q-path exists

Containment for XPath(/, //) and DTDs _

Result 1 [Neven, Sch. 2003]

The containment problem for XPath(/, //) in the presence of DTDs is in **PTIME**

Proof idea (cont.)

- Deterministic top-down automaton \mathcal{A}_p
- Deterministic top-down automaton $\mathcal{A}_{\overline{q}}$ checking that no q-path exists
- There is a deterministic top-down automaton \mathcal{A}_d checking whether t conforms to d
- $p \subseteq_b q$ in the presence of $d \Longleftrightarrow L(\mathcal{A}_p \times \mathcal{A}_{\overline{q}} \times \mathcal{A}_d) = \emptyset$
- The latter can be checked in polynomial time



Schwentick

Containment for XPath(/,//,[],*,|) and DTDs ____

Lemma For each XPath(/,//,[], *, |)-expression p there is a deterministic bottom-up automaton \mathcal{A}_p of exponential size which checks whether in a tree p holds

Proof idea for Lemma

- States of \mathcal{A}_p are of the form $(S_{/},S_{//})$
- Both $S_{/}$ and $S_{//}$ are sets of positions of the query tree:
 - $-S_{\prime}$: positions matching v
 - $S_{//}$: positions matching some node in the subtree of v

Containment for XPath(/,//,[],*,|) and DTDs ____

Result 2 [Neven, Sch. 2003]

The containment problem for XPath(/, //, [], *, |) in the presence of DTDs is in **EXPTIME**

Proof idea (cont.)

- Construct deterministic bottom-up automaton \mathcal{A}_p of exponential size
- Construct deterministic bottom-up automaton $\mathcal{A}_{\overline{q}}$ of exponential size
- Construct deterministic bottom-up automaton \mathcal{A}_d of exponential size
- $p \subseteq_b q$ in the presence of $d \Longleftrightarrow L(\mathcal{A}_p \times \mathcal{A}_{\overline{q}} \times \mathcal{A}_d) = \emptyset$

 \Rightarrow exponential time















Strategies as trees ____



Strategies as trees





Winning strategies and paths

Proof Sketch (cont.)

There are various kinds of paths in a game tree:

- (a) Legal tilings \implies Player I wins
- (b) Syntactically wrong: some row of wrong length
- (c) II places a wrong tile \implies Player I wins
- (d) I places a wrong tile \implies Player II wins

Player I has a winning strategy

\iff

there is a tree in which all paths are of the form (a) or (c)

We want to construct q such that all paths of the form (b) or (d) are selected Then: Player I wins iff $/S \not\subseteq q$ wrt DTD

Winning strategies and paths

Proof Sketch (cont.)

There are various kinds of paths in a game tree:

- (a) Legal tilings \implies Player I wins
- (b) Syntactically wrong: some row of wrong length
- (c) II places a wrong tile \implies Player I wins
- (d) I places a wrong tile \implies Player II wins

Player I has a winning strategy

\iff

there is a tree in which all paths are of the form (a) or (c)

We want to construct q such that all paths of the form (b) or (d) are selected Then: Player I wins iff $/S \not\subseteq q$ wrt DTD

Problem: if II places a wrong tile, I might be forced to place a wrong tile, too

 \implies We let player I mark wrong tiles of II by !

 \implies We have to check that I does this correctly

Proof Sketch (cont.)

Player I has winning strategy $\iff /S \not\subseteq q$

q expresses that one of the following holds

- Player I violates a horizontal constraint:
- Player I violates a vertical constraint:
- Some row does not contain exactly n tiles
- Player I wrongly claims a mistake of II:
- Some more conditions on ${\boldsymbol{B}}$ and ${\boldsymbol{T}}$

Proof Sketch (cont.)

Player I has winning strategy $\iff /S \not\subseteq q$

q expresses that one of the following holds

• Player I violates a horizontal constraint:

For each $(x,y) \not\in H$: //(x,II)/(y,I)

- Player I violates a vertical constraint:
- Some row does not contain exactly n tiles
- Player I wrongly claims a mistake of II:
- Some more conditions on $m{B}$ and $m{T}$
- *= OR of all symbols, $\sigma^i=\sigma/\cdots/\sigma$ (i times)

Proof Sketch (cont.)

Player I has winning strategy $\iff /S \not\subseteq q$

 \boldsymbol{q} expresses that one of the following holds

• Player I violates a horizontal constraint:

For each $(x,y)
ot\in H$: //(x,II)/(y,I)

• Player I violates a vertical constraint:

For each $(x,y) \not\in V$: $//(x,I)/*^{n+1}/(y,I)$

- Some row does not contain exactly $m{n}$ tiles
- Player I wrongly claims a mistake of II:
- Some more conditions on **B** and **T**
- *= OR of all symbols, $\sigma^i=\sigma/\cdots/\sigma$ (i times)

Proof Sketch (cont.)

Player I has winning strategy $\iff /S \not\subseteq q$

 \boldsymbol{q} expresses that one of the following holds

• Player I violates a horizontal constraint:

For each $(x,y)
ot\in H$: //(x,II)/(y,I)

• Player I violates a vertical constraint:

For each $(x,y) \not\in V$: $//(x,I)/*^{n+1}/(y,I)$

• Some row does not contain exactly n tiles

$${\mathcal D}^{n+1} \mid igcup_{i=0}^{n-1}(\$^{\scriptscriptstyle {||}} \$^{\scriptscriptstyle {||}} |S)/{\mathcal D}^{i}/(\$^{\scriptscriptstyle {||}} \$^{\scriptscriptstyle {||}} |\#)$$

- Player I wrongly claims a mistake of II:
- Some more conditions on B and T

$$*=$$
 OR of all symbols, $\sigma^i=\sigma/\cdots/\sigma$ (i times) $\mathcal{D}=(d_1,I)|\cdots|(d_m,I)|(d_1,II)|\cdots|(d_m,II)$

Proof Sketch (cont.)

Player I has winning strategy $\iff /S \not\subseteq q$

 \boldsymbol{q} expresses that one of the following holds

• Player I violates a horizontal constraint:

For each $(x,y)
ot\in H$: //(x,II)/(y,I)

• Player I violates a vertical constraint:

For each $(x,y) \not\in V$: $//(x,I)/*^{n+1}/(y,I)$

• Some row does not contain exactly $m{n}$ tiles

$${\mathcal D}^{n+1} \mid igcup_{i=0}^{n-1}(\$^{ ext{!`}}|\$^{ ext{!`'}}|S)/{\mathcal D}^{i}/(\$^{ ext{!`}}|\$^{ ext{!`'}}|\#)$$

• Player I wrongly claims a mistake of II:

For each $(x,y) \in V, (x',y) \in H$: $//(x,II)/*^n/(x',I)/(y,II)/!$

• Some more conditions on ${\pmb B}$ and ${\pmb T}$

$$*=$$
 OR of all symbols, $\sigma^i=\sigma/\cdots/\sigma$ (i times)

 $= (d_1,I)|\cdots |(d_m,I)|(d_1,II)|\cdots |(d_m,II)|$

Schwentick

Related work on XPath containment

More Results

- Containment of XPath with / and a subset of {//,[],*} was studied in [Miklau and Suciu 2002]:
 - Containment of XPath(//,[],*) is coNP-complete even if the number of
 * or the number of [] is bounded
 - If the number of // is bounded then it is in polynomial time
- XPath containment in the presence of DTDs and simple integrity constraints was investigated in [Deutsch and Tanen 2001]:
 - In general (unbounded constraints): undecidable
- More complexity results between **coNP** and undecidable for other fragments and extensions in [Neven and S. 2003]

Some Open Questions

- What's the exact borderline between fragments of XPath with decidable and undecidable containment problem?
- To what extent can the presented result be extended to other axes (siblings, backward)?
Summary: XPath

Summary

Expressive Power

Closely related to first-order logic

Evaluation

- In general: Polynomial time
- Large fragments in linear time
- Structural complexity between LOGSPACE and PTIME

Satisfiability

- Without negation: **PTIME** or **NP**
- With negation: non-elementary

Containment

- Varying from **PTIME** to undecidable
- Upper bound for positive navigational XPath? XML: Algorithms & Complexity
 Introduction - XPath

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

XSLT Typechecking

Definition: Transformation typechecking

Given: DTDs d_1 and d_2 and a transformation T

Result: Is T(t) valid wrt. d_2 , for each document t valid wrt. d_1 ?

XSLT Typechecking

Definition: Transformation typechecking

Given: DTDs d_1 and d_2 and a transformation T

Result: Is T(t) valid wrt. d_2 , for each document t valid wrt. d_1 ?

Question: Is XSLT typechecking decidable?

Question: What is the complexity?

XSLT Typechecking

Definition: Transformation typechecking

- **Given:** DTDs d_1 and d_2 and a transformation T
- **Result:** Is T(t) valid wrt. d_2 , for each document t valid wrt. d_1 ?

Question: Is XSLT typechecking decidable?

Question: What is the complexity?

Outline of the Following

- Provide an automata model for XSLT transformations
- Show that the behaviour of these automata can be captured by MSO logic
- Use manipulation of regular tree languages to solve type checking problem
- → This part is based on [Milo,Suciu,Vianu 01]

XSLT in more detail

How XSLT Roughly Works

Templates:

(xsl:template name=TName match=pattern mode=MName)

Template application:

{xsl:apply-templates select=Expression mode=MName}

XSLT Processing Whenever xsl:apply-templates is called at a node v the following happens:

- Compute set S(v) of nodes, reachable from v via Expression (if select is not present, S(v) = children of v)
- For each $w \in S(v)$ compute which templates that can be applied to w:
 - \boldsymbol{w} has to match pattern of a template
 - the mode of the template has to be the same as the mode of xsl:apply-templates
- If no template matches, take the default template
- For each $w \in S(v)$ select the best template and apply it.

The process starts at the root of the tree

Example Transformation

Remove everything below a c. Translate a below b into d

Example XSLT

```
\langle xsl:template match="a" \rangle
    \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle
\langle xsl:template \rangle
(xsl:template match="a" mode="below")
    \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle
(/xsl:template)
\langle xsl:template match="b" \rangle
    \langle b \rangle \langle xsl:apply-templates mode="below" \rangle \langle b \rangle
\langle xsl:template \rangle
(xsl:template match="b" mode="below")
    \langle b \rangle \langle xsl:apply-templates mode="below" \rangle \langle b \rangle
\langle xsl:template \rangle
\langle xsl:template match="c" \rangle
    \langle c \rangle \langle /c \rangle
\langle xsl:template \rangle
\langle xsl:template match="c" mode="below" \rangle
    \langle c \rangle \langle /c \rangle
(/xsl:template)
```

Schwentick

XML: Algorithms & Complexity

Example Transformation

Remove everything below a c. Translate a below b into d

Example XSLT (Abbreviated)

 $\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$

Example Transformation

Remove everything below a c. Translate a below b into d

Example XSLT (Abbreviated)

 $\begin{array}{l} & \left(... \mbox{ match}{=}"a" \right) \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ & \left(... \mbox{ match}{=}"a" \mbox{ mode}{=}"below" \right) \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ & \left(... \mbox{ match}{=}"b" \right) \langle b \rangle \langle xsl:apply-templates \mbox{ mode}{=}"below" \rangle \langle /b \rangle \langle /... \rangle \\ & \left(... \mbox{ match}{=}"b" \mbox{ mod}{=}"below" \right) \langle b \rangle \langle xsl:apply-templates \mbox{ mod}{=}"below" \rangle \langle /b \rangle \langle /... \rangle \\ & \left(... \mbox{ match}{=}"c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ & \left(... \mbox{ match}{=}"c" \mbox{ mod}{=}"below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \end{array} \right)$

Example Trees

Example Transformation

Remove everything below a c. Translate a below b into d

$$\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$$



Example Transformation

Remove everything below a c. Translate a below b into d

$$\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$$



Example Transformation

Remove everything below a c. Translate a below b into d

$$\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$$



Example Transformation

Remove everything below a c. Translate a below b into d

$$\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$$



Example Transformation

Remove everything below a c. Translate a below b into d

$$\begin{array}{l} \langle ... \mbox{ match}="a" \rangle \langle a \rangle \langle xsl:apply-templates \rangle \langle /a \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="a" \mbox{ mode}="below" \rangle \langle d \rangle \langle xsl:apply-templates \rangle \langle /d \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="b" \mbox{ mode}="below" \rangle \langle b \rangle \langle xsl:apply-templates \mbox{ mode}="below" \rangle \langle /b \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle \\ \langle ... \mbox{ match}="c" \mbox{ mode}="below" \rangle \langle c \rangle \langle /c \rangle \langle /... \rangle$$





Example Trees



Example Trees























An automaton model for XSLT

Definition: *k*-pebble Transducer

- Work on binary tree encodings of unranked trees
- Up to k pebbles can be placed on the tree
- Only pebble with highest number (*current pebble*) can move, depending on state, number of pebbles symbols under pebbles and incidence of pebbles
- possible pebble movements:
 - stay
 - go to left child, right child or parent
 - lift current pebble
 - place new pebble on the root
- Nondeterminism allowed
- If current pebble stays it is possible to produce output:
 - a node with two (forthcoming) subtrees; in this case two independent subcomputations (*branches*) are started, which construct the left subtree and right subtree, respectively
 - a leaf; in this case the computation branch stops

An automaton model for XSLT



An automaton model for XSLT

Definition: *k*-pebble Transducer

- Work on binary tree encodings of unranked trees
- Up to k pebbles can be placed on the tree
- Only pebble with highest number (*current pebble*) can move, depending on state, number of pebbles symbols under pebbles and incidence of pebbles
- possible pebble movements:
 - stay
 - go to left child, right child or parent
 - lift current pebble
 - place new pebble on the root
- Nondeterminism allowed
- If current pebble stays it is possible to produce output:
 - a node with two (forthcoming) subtrees; in this case two independent subcomputations (*branches*) are started, which construct the left subtree and right subtree, respectively
 - a leaf; in this case the computation branch stops

Computing XSLT transformations by *k***-pebble transducers**

Fact

k-pebble transducers can evaluate most XPath expressions

(and produce as output an encoded version of the result list)

- even with other axes than the forward axis

- Whenever xsl:apply-templates is called at a node v the following happens:
 - Cycle through the set S(v) of nodes, reachable from v via Expression (if select is not present, S(v) = children of v)
 - For each $w \in S(v)$ check which templates can be applied to w:
 - $\ast \ oldsymbol{w}$ has to match pattern of a template
 - * the mode of xsl:apply-templates is stored in the state of the automaton
 - For each $w \in S(v)$ select the best template and branch into
 - * a subcomputation which handles the next node in S(v) (via the right child)
 - * a subcomputation which applies the template to the current node
- The computation starts at the root of the tree

Back to the Typechecking Question

Question: Is XSLT typechecking decidable?

- How can we check that $T(t) \in L(d_2)$, for each $t \in L(d_1)$?
- Obvious approach:
 - Compute $T(L(d_1))$
 - Check that $T(L(d_1)) \subseteq L(d_2)$

Back to the Typechecking Question

Question: Is XSLT typechecking decidable?

- How can we check that $T(t) \in L(d_2)$, for each $t \in L(d_1)$?
- Obvious approach:
 - Compute $T(L(d_1))$
 - Check that $T(L(d_1)) \subseteq L(d_2)$
- Problem: $T(L(d_1))$ does not need to be regular:



Back to the Typechecking Question

Question: Is XSLT typechecking decidable?

- How can we check that $T(t) \in L(d_2)$, for each $t \in L(d_1)$?
- Obvious approach:
 - Compute $T(L(d_1))$
 - Check that $T(L(d_1)) \subseteq L(d_2)$
- Problem: $T(L(d_1))$ does not need to be regular:



- Alternative approach:
 - Compute $T^{-1}(\overline{L(d_2)})$
 - Check $L(d_1) \cap T^{-1}(\overline{L(d_2)}) = \emptyset$



- Basically the same as *k*-pebble transducers
- Instead of output producing steps:
 - accept
 - branch into two independent subcomputations
- A tree is accepted if all subcomputations accept

Main Steps of the Proof

(i) $T^{-1}(\overline{L(d_2)})$ is accepted by a *k*-pebble acceptor

(ii) k-pebble acceptors only accept regular tree languages

Step (i)

Lemma

 $T^{-1}(\overline{L(d_2)})$ is accepted by a *k*-pebble acceptor

Proof

- Let B be a nondeterministic top-down tree automaton which accepts $\overline{L(d_2)}$
- Let T be a k-pebble tree transducer
- We construct k-pebble acceptor A for $T^{-1}(\overline{L(d_2)})$, i.e., an automaton which on input t decides whether there is a tree in T(t) which is accepted by B:
 - Simulate T on t and B
 - Simulate at the same time the behaviour of *B* on the (virtual) output tree
 this is possible as the output tree is produced top-down and can be instantly consumed by *B*
 - The simulation involves branching, whenever T branches, and produces two new subtrees

Step (ii) ____

Lemma
$m{k}$ -pebble acceptors only accept regular tree languages
Proof idea
Show that the language of a $m k$ -pebble acceptor can be
expressed by an MSO-formula:
1. Reduce k -pebble automaton acceptance to AGAP
(Alternating Graph Accessibility)

- 2. Show that AGAP can be expressed in MSO
- 3. Some adjustments necessary

Definition: Accessible Nodes

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Definition: Accessible Nodes

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Definition: Accessible Nodes

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Definition: Accessible Nodes

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Definition: Accessible Nodes

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible


Alternating Graph Accessibility

Definition: Accessible Nodes

Let G = (V, E), $V = V_{\wedge} \cup V_{\vee}$. A node w is accesible if

- ullet $w \in V_{\wedge}$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Alternating Graph Accessibility

Definition: Accessible Nodes

Let G = (V, E), $V = V_{\wedge} \cup V_{\vee}$. A node w is accesible if

- ullet $w \in V_\wedge$ and all successors of w are accessible, or
- ullet $w \in V_{ee}$ and at least one successor of w is accessible



Alternating Graph Accessibility (cont.)

Construction of $G_{A,t}$ from Automaton A and Tree t

- Nodes in V_{\vee} are the configurations of A on t: tuples $[i, q, \theta]$, where $\theta : \{1, \dots, i\} \to t$
- Nodes in V_{\wedge} are ϵ and pairs (γ_1, γ_2) of configurations with "the same heta"

• Edges:
$$-(\gamma_1, \gamma_2) \rightarrow \gamma_1$$
, $(\gamma_1, \gamma_2) \rightarrow \gamma_2$

– $\gamma
ightarrow \gamma'$, if this is a step of A

– $\gamma
ightarrow \epsilon$, if A can get into the accept state from γ

 $-\gamma
ightarrow (\gamma_1,\gamma_2)$ if this is a branching step of A

Fact

A k-pebble acceptor A accepts a tree $t \Longleftrightarrow \gamma$ is accessible in $(G_{A,t})$

AGAP is MSO-expressible

Definition: Reverse-closed Sets of Nodes

A set S of nodes is reverse-closed if the following holds

- ullet if v is in V_\wedge and $w\in S$, for all nodes w with $(v,w)\in E$, then $v\in S$
- ullet if v is in V_{ee} and $w\in S$, for some node w with $(v,w)\in E$, then $v\in S$



Fact
Node $oldsymbol{v}$ is accessible iff it is in every reverse-closed
set of nodes.

...as MSO-Formula

v accessible $\equiv orall S$ (reverse-closed(S)
ightarrow S(v)), where

 $egin{aligned} \mathsf{reverse-closed}(S) &\equiv & orall x \left([V_\wedge(x) \wedge orall y \left(E(x,y) o S(y)
ight)] o S(x)
ight) \wedge \ & \left([V_ee(x) \wedge \exists y \left(E(x,y) \wedge S(y)
ight)] o S(x)
ight) \end{aligned}$

k-pebble acceptors and MSO

Proof idea

Unfortunately, $G_{A,t}$ has too many nodes to use this directly:

- MSO can only quantify over sets of linear size in the given structure (i.e., t)
- $G_{A,t}$ has $\Omega(|t|^k)$ configurations
- But G_{A,t} has a special structure: Nodes are only connected if their number of pebbles is the same ±1 and they agree in all but at most the last pebble

k-pebble acceptors and MSO (cont.)

Proof (cont.)

- Wlog assume that each state of A is only used for a fixed number of pebbles: $Q = Q_1 \cup \cdots \cup Q_k$, where the states in Q_i are only used, when i pebbles are present
- Further assume that all sets Q_i are of equal size m: $Q_i = \{q_{i1}, \ldots, q_{im}\}$

• k = 1:

- Use one relation S_i^1 for each state q_{1i}
- Intended meaning of $v \in S_i^1$: there is an accepting subcomputation of A starting at v in state q_{1i}
- $\varphi = \forall S_1^1 \cdots \forall S_m^1 \; (\text{reverse-closed} \to S_1^1(\text{root}))$
- reverse-closed is a conjunction of subformulas, induced by the transitions of *A*, e.g.:
 - $* \hspace{0.1 cm} ext{if} \hspace{0.1 cm} (q_{1i},a) o ext{accept then} \hspace{0.1 cm} orall x \hspace{0.1 cm} Q_a(x) o S_i^1(x)$
 - $* \hspace{0.1 in} ext{if} \hspace{0.1 in} (q_{1i}, a)
 ightarrow (q_{1j}, ext{down-right}) ext{ then}$

 $\forall x \,\forall y \, (Q_a(x) \wedge E_r(x,y) \wedge S_j^1(y)) \to S_i^1(x)$

k-pebble acceptors and MSO (cont.)

Proof (cont.)

k = 2:

- reverse-closed¹ and reverse-closed² describe reverse closure for configurations with one and two pebbles, respectively
- reverse-closed² expresses the same as reverse-closed before, but with the (immobile) pebble 1 represented by variable x_1
- reverse-closed¹ also refers to subcomputations with a second pebble
- Conjuncts corresponding to simple movements are essentially the same
- Conditions which check whether pebbles are at the same node have to be added
- The following conjuncts are added for pebble placement and lifting:
 - $(q_{2i},a) \to (q_{1j},{\rm lift})$ adds $orall x_2 \left(Q_a(x_2) \wedge S_j^1(x_1)\right) \to S_i^2(x_2)$ to reverse-closed 2
 - $(q_{1i}, a) \rightarrow (q_{2j}, \text{place})$ adds $\forall x_1 (Q_a(x_1) \land \varphi^2) \rightarrow S_i^1(x_1)$ to reverse-closed¹, where φ^2 is $\forall S_1^2 \cdots \forall S_m^2$ (reverse-closed² $\rightarrow S_j^2$ (root))

Summary of Proof

Proof (cont.)
$ullet$ To solve the type checking problem, given d_1 , d_2 and T , we can proceed as
follows.
(1) Construct the k -pebble acceptor A for $T^{-1}(\overline{L(d_2)})$
(2) Transform A into an equivalent MSO formula Φ
(3) Φ holds for all trees t for which $T(t) ot \subseteq L(d_2)$
(4) Construct a nondeterministic bottom-up automaton A' equivalent to $ eg \Phi$
(5) Check that $L(d_1) \subseteq L(A')$
 Hence, the type-checking problem is decidable
• Steps (1) and (4) can be done in poly-time

- Step (2) is exponential in k, FO-quantifier depth of Φ is k, MSO-quantifier depth of Φ is |Q|
- Step (3) is non-elementary (exponentiation tower of height k)
- Hence,

the algorithm for the type-checking problem has a very bad complexity

Summary: Typechecking

Related Work

- If transformations are allowed to compare data values in the input document type checking becomes undecidable very quickly, even for restricted types and transformations [Alon et al. 2001]
- Typechecking for deterministic top-down tree transducers is more tractable. Complexity depends on exact representation of DTDs and restrictions on the transducers: between **PTIME** and **EXPTIME** [Martens and Neven 2003]
- If $P \neq NP$ there is no elementary f, such that MSO-formulas can be evaluated in time $f(|formula|) \times p(|tree|)$ with polynomial p [Frick and Grohe 2002]

Open

- Find (more) transformations with a tractable typechecking problem
- In particular, with data values

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Theory of XQuery

Remarks

- In general, the theoretical foundations of XQuery have to be developed
- Clearly: XQuery is Turing-complete and therefore static analysis is generally impossible
- What about important fragments with better properties?
- E.g., Tree pattern queries
- Here, we concentrate on:
 - Conjunctive queries for trees
 - Some questions related to automata for XQuery

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conjunctive Queries

Automata and XQuery

Conclusion

Conjunctive Queries

Introduction

- Navigational XPath expressions (without or and not) can be written as conjunctive queries
- /child::a/desc::*[child::c]/parent::d corresponds to $Q(x) = \operatorname{root}(x_1) \wedge \operatorname{child}(x_1, x_2) \wedge L_a(x_2) \wedge \operatorname{desc}(x_2, x_3) \wedge \operatorname{child}(x_3, x_4) \wedge L_c(x_4) \wedge \operatorname{child}(x, x_3) \wedge L_d(x)$
- Conjunctive Queries can express queries of higher arity: $Q(x,y) = ext{child}(x,x_1) \wedge ext{child}(x_1,y)$
- What is the complexity of evaluating conjunctive queries on trees?
- Data complexity is in **PTIME** (even in **LOGSPACE**): Cycle through all valuations of the variables
- What about combined complexity?

A Generic Algorithm

Definition

- Pre-valuation : mapping from variables to non-empty sets of nodes
- For a conjunctive query Q a pre-valuation θ is consistent if:
 - for each atom $L_{\sigma}(x)$: $v \in heta(x) \Rightarrow L_{\sigma}(v)$
 - for each atom R(x, y):
 - $* \ v \in heta(x) \Rightarrow \exists u \in heta(y) R(v,u)$
 - $* \ v \in heta(y) \Rightarrow \exists u \in heta(x) R(u,v)$





A Generic Algorithm

Definition

- Pre-valuation : mapping from variables to non-empty sets of nodes
- For a conjunctive query Q a pre-valuation θ is consistent if:
 - for each atom $L_{\sigma}(x)$: $v \in heta(x) \Rightarrow L_{\sigma}(v)$
 - for each atom R(x, y):
 - $* \ v \in heta(x) \Rightarrow \exists u \in heta(y) R(v,u)$
 - $v \in heta(y) \Rightarrow \exists u \in heta(x) R(u,v)$

















Question: Is *h* always a solution?





Definition

A binary relation R is *<*-hemichordal if for all u, u', v, v' with u < u' and $u \le v \le v'$

- ullet $R(u,v') \wedge R(u',v)
 ightarrow R(u,v)$ and
- ullet $R(v',u) \wedge R(v,u')
 ightarrow R(v,u)$

Theorem [Gottlob, Koch, Schulz 04]

If the relations of a query Q are <-hemichordal and θ is a

consistent pre-valuation for Q

then the <-minimal valuation for θ is a solution for Q

Definition A binary relation R is *<*-hemichordal if for all u, u', v, v' with u < u' and $u \le v \le v'$

- ullet $R(u,v') \wedge R(u',v)
 ightarrow R(u,v)$ and
- ullet $R(v',u) \wedge R(v,u')
 ightarrow R(v,u)$



Combined Complexity of Conjunctive Queries

Observation

It is sufficient to consider the axes child, child⁺, child^{*}, NextSibling, NextSibling⁺, NextSibling^{*}, Following

Theorem [Gottlob, Koch, Schulz 04]

- child⁺ and child^{*} are preorder-hemichordal
- following is postorder-hemichordal
- child, NextSibling, NextSibling⁺, NextSibling^{*} are breadth-first-left-to-right-hemichordal

Corollary

For each of these sets of axes conjunctive queries can be evaluated in time $O(query size \times tree size)$

Combined Complexity of Conjunctive Queries

Observation

It is sufficient to consider the axes child, child⁺, child^{*}, NextSibling, NextSibling⁺, NextSibling^{*}, Following

Theorem [Gottlob, Koch, Schulz 04]

- child⁺ and child^{*} are preorder-hemichordal
- following is postorder-hemichordal
- child, NextSibling, NextSibling⁺, NextSibling^{*} are breadth-first-left-to-right-hemichordal

Corollary

For each of these sets of axes conjunctive queries can be evaluated in time $O(query size \times tree size)$

Amazing Result

For sets of axes not contained in those, the combined complexity of conjunctive query evaluation is **NP**-complete

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conjunctive Queries

Automata and XQuery

Conclusion

Automata and XQuery _

So far...

- We have seen that automata are useful for
 - Validation, Typing
 - Navigation
 - Transformation
- What about more general queries?
 - results of higher arity?
 - joins, i.e., comparisons of data values
 - counting
- Are automata useful for XQuery?
- ... for tree pattern queries?

General Queries (cont.)

Higher arity

- Nonemptiness and containment questions can be handled by automata: tuples can be encoded by additional labels
- What about query evaluation for higher arity?

Data values

- When data values in XML documents are taken into account, things become more complicated, e.g.:
 - Even First-order logic becomes undecidable
 - Pebble automata become undecidable
 - Automata with data registers become undecidable when they are allowed to move up and down
- What is the right notion for regular (string) languages over infinite alphabets?
- What are sensible decidable restrictions of logics and automata in the context of data values?

General Queries (cont.)

Counting

- Automata can be equipped with counting facilities, e.g.: Presburger tree automata: $\delta(\sigma, q)$ is Boolean combination of
 - regular expressions and
 - quantifier-free Presburger formulas like

"number of children in state q_1 = number of children in state q_2 "

- Nondet. Presburger automata:
 - $\equiv MSO \log ic$
 - Whether automaton accepts all trees is undecidable
- Det. Presburger automata:
 - \equiv Presburger μ -formulas
 - Membership test: $O(|\mathcal{A}||t|)$
 - Non-emptiness: **PSPACE**
 - Containment: **PSPACE**

[Seidl, Sch., Muscholl, Habermehl 2004]

Contents

Introduction

Background on Tree Automata and Logic

Schema Languages

XPath and Node-selecting Queries

XSLT

XQuery

Conclusion

Summary

- Schema languages and XPath are well understood
- There are some nice results on transformations
- Theory for XQuery still has to be developed

Summary

- Schema languages and XPath are well understood
- There are some nice results on transformations
- Theory for XQuery still has to be developed

Finally...

Thanks for your patience